

how it works, but you don't have to code this part yourself—it underlies the R or Python functions.

Overfitting

Throughout the book you will be cautioned repeatedly about *overfitting*, possibly to the point you will have nightmares about it. Overfitting is the term used to mean that you used a dataset to estimate the parameters of your model, but your model isn't that good at capturing reality beyond your sampled data.

You might know this because you have tried to use it to predict labels for another set of data that you didn't use to fit the model, and it doesn't do a good job, as measured by an evaluation metric such as accuracy.

Exploratory Data Analysis

“Exploratory data analysis” is an attitude, a state of flexibility, a willingness to look for those things that we believe are not there, as well as those we believe to be there.

— John Tukey

Earlier we mentioned exploratory data analysis (EDA) as the first step toward building a model. EDA is often relegated to chapter 1 (by which we mean the “easiest” and lowest level) of standard introductory statistics textbooks and then forgotten about for the rest of the book.

It's traditionally presented as a bunch of histograms and stem-and-leaf plots. They teach that stuff to kids in fifth grade so it seems trivial, right? No wonder no one thinks much of it.

But EDA is a critical part of the data science process, and also represents a philosophy or way of doing statistics practiced by a strain of statisticians coming from the Bell Labs tradition.

John Tukey, a mathematician at Bell Labs, developed exploratory data analysis in contrast to confirmatory data analysis, which concerns itself with modeling and hypotheses as described in the previous section. In EDA, there is no hypothesis and there is no model. The “exploratory” aspect means that your understanding of the problem you are solving, or might solve, is changing as you go.

Historical Perspective: Bell Labs

Bell Labs is a research lab going back to the 1920s that has made innovations in physics, computer science, statistics, and math, producing languages like C++, and many Nobel Prize winners as well. There was a very successful and productive statistics group there, and among its many notable members was John Tukey, a mathematician who worked on a lot of statistical problems. He is considered the father of EDA and R (which started as the S language at Bell Labs; R is the open source version), and he was interested in trying to visualize high-dimensional data.

We think of Bell Labs as one of the places where data science was “born” because of the collaboration between disciplines, and the massive amounts of complex data available to people working there. It was a virtual playground for statisticians and computer scientists, much like Google is today.

In fact, in 2001, Bill Cleveland wrote “Data Science: An Action Plan for expanding the technical areas of the field of statistics,” which described multidisciplinary investigation, models, and methods for data (traditional applied stats), computing with data (hardware, software, algorithms, coding), pedagogy, tool evaluation (staying on top of current trends in technology), and theory (the math behind the data).

You can read more about Bell Labs in the book *The Idea Factory* by Jon Gertner (Penguin Books).

The basic tools of EDA are plots, graphs and summary statistics. Generally speaking, it’s a method of systematically going through the data, plotting distributions of all variables (using box plots), plotting time series of data, transforming variables, looking at all pairwise relationships between variables using scatterplot matrices, and generating summary statistics for all of them. At the very least that would mean computing their mean, minimum, maximum, the upper and lower quartiles, and identifying outliers.

But as much as EDA is a set of tools, it’s also a mindset. And that mindset is about your relationship with the data. You want to understand the data—gain intuition, understand the shape of it, and try to connect your understanding of the process that generated the data to

the data itself. EDA happens between you and the data and isn't about proving anything to anyone else yet.

Philosophy of Exploratory Data Analysis

Long before worrying about how to convince others, you first have to understand what's happening yourself.

— Andrew Gelman

While at Google, Rachel was fortunate to work alongside two former Bell Labs/AT&T statisticians—Daryl Pregibon and Diane Lambert, who also work in this vein of applied statistics—and learned from them to make EDA a part of her best practices.

Yes, even with very large Google-scale data, they did EDA. In the context of data in an Internet/engineering company, EDA is done for some of the same reasons it's done with smaller datasets, but there are additional reasons to do it with data that has been generated from logs.

There are important reasons anyone working with data should do EDA. Namely, to gain intuition about the data; to make comparisons between distributions; for sanity checking (making sure the data is on the scale you expect, in the format you thought it should be); to find out where data is missing or if there are outliers; and to summarize the data.

In the context of data generated from logs, EDA also helps with debugging the logging process. For example, “patterns” you find in the data could actually be something wrong in the logging process that needs to be fixed. If you never go to the trouble of debugging, you'll continue to think your patterns are real. The engineers we've worked with are always grateful for help in this area.

In the end, EDA helps you make sure the product is performing as intended.

Although there's lots of visualization involved in EDA, we distinguish between EDA and data visualization in that EDA is done toward the beginning of analysis, and data visualization (which we'll get to in [Chapter 9](#)), as it's used in our vernacular, is done toward the end to communicate one's findings. With EDA, the graphics are solely done for *you* to understand what's going on.

With EDA, you can also use the understanding you get to inform and improve the development of algorithms. For example, suppose you

are trying to develop a ranking algorithm that ranks content that you are showing to users. To do this you might want to develop a notion of “popular.”

Before you decide how to quantify popularity (which could be, for example, highest frequency of clicks, or the post with the most number of comments, or comments above some threshold, or some weighted average of many metrics), you need to understand how the data is behaving, and the best way to do that is looking at it and getting your hands dirty.

Plotting data and making comparisons can get you extremely far, and is far better to do than getting a dataset and immediately running a regression just because you know how. It’s been a disservice to analysts and data scientists that EDA has not been enforced as a critical part of the process of working with data. Take this opportunity to make it part of your process!

Here are some references to help you understand best practices and historical context:

1. *Exploratory Data Analysis* by John Tukey (Pearson)
2. *The Visual Display of Quantitative Information* by Edward Tufte (Graphics Press)
3. *The Elements of Graphing Data* by William S. Cleveland (Hobart Press)
4. *Statistical Graphics for Visualizing Multivariate Data* by William G. Jacoby (Sage)
5. “Exploratory Data Analysis for Complex Models” by Andrew Gelman (American Statistical Association)
6. *The Future of Data Analysis* by John Tukey. *Annals of Mathematical Statistics*, Volume 33, Number 1 (1962), 1-67.
7. *Data Analysis, Exploratory* by David Brillinger [8-page excerpt from *International Encyclopedia of Political Science* (Sage)]

Exercise: EDA

There are 31 datasets named `nyt1.csv`, `nyt2.csv`, ..., `nyt31.csv`, which you can find here: https://github.com/oreillymedia/doing_data_science.

Each one represents one (simulated) day's worth of ads shown and clicks recorded on the *New York Times* home page in May 2012. Each row represents a single user. There are five columns: age, gender (0=female, 1=male), number impressions, number clicks, and logged-in.

You'll be using R to handle these data. It's a programming language designed specifically for data analysis, and it's pretty intuitive to start using. You can download it [here](#). Once you have it installed, you can load a single file into R with this command:

```
data1 <- read.csv(url("http://stat.columbia.edu/~rachel/datasets/nyt1.csv"))
```

Once you have the data loaded, it's time for some EDA:

1. Create a new variable, `age_group`, that categorizes users as "<18", "18-24", "25-34", "35-44", "45-54", "55-64", and "65+".
2. For a single day:
 - Plot the distributions of number impressions and click-through-rate (CTR=# clicks/# impressions) for these six age categories.
 - Define a new variable to segment or categorize users based on their click behavior.
 - Explore the data and make visual and quantitative comparisons across user segments/demographics (<18-year-old males versus < 18-year-old females or logged-in versus not, for example).
 - Create metrics/measurements/statistics that summarize the data. Examples of potential metrics include CTR, quantiles, mean, median, variance, and max, and these can be calculated across the various user segments. Be selective. Think about what will be important to track over time—what will compress the data, but still capture user behavior.
3. Now extend your analysis across days. Visualize some metrics and distributions over time.
4. Describe and interpret any patterns you find.

Sample code

Here we'll give you the beginning of a sample solution for this exercise. The reality is that we can't teach you about data science and teach you

how to code all in the same book. Learning to code in a new language requires a lot of trial and error as well as going online and searching on Google or stackoverflow.

Chances are, if you're trying to figure out how to plot something or build a model in R, other people have tried as well, so rather than banging your head against the wall, look online. [Ed note: There might also be some **books** available to help you out on this front as well.] We suggest not looking at this code until you've struggled along a bit:

```
# Author: Maura Fitzgerald
data1 <- read.csv(url("http://stat.columbia.edu/~rachel/
                    datasets/nyt1.csv"))

# categorize
head(data1)
data1$agecat <- cut(data1$Age, c(-Inf, 0, 18, 24, 34, 44, 54, 64, Inf))

# view
summary(data1)

# brackets
install.packages("doBy")
library("doBy")
siterange <- function(x){c(length(x), min(x), mean(x), max(x))}
summaryBy(Age~agecat, data =data1, FUN=siterange)

# so only signed in users have ages and genders
summaryBy(Gender+Signed_In+Impressions+Clicks~agecat,
          data =data1)

# plot
install.packages("ggplot2")
library(ggplot2)
ggplot(data1, aes(x=Impressions, fill=agecat))
  +geom_histogram(binwidth=1)
ggplot(data1, aes(x=agecat, y=Impressions, fill=agecat))
  +geom_boxplot()

# create click thru rate
# we don't care about clicks if there are no impressions
# if there are clicks with noimps my assumptions about
# this data are wrong
data1$hasimps <- cut(data1$Impressions, c(-Inf, 0, Inf))
summaryBy(Clicks~hasimps, data =data1, FUN=siterange)
ggplot(subset(data1, Impressions>0), aes(x=Clicks/Impressions,
    colour=agecat)) + geom_density()
ggplot(subset(data1, Clicks>0), aes(x=Clicks/Impressions,
    colour=agecat)) + geom_density()
ggplot(subset(data1, Clicks>0), aes(x=agecat, y=Clicks,
```

```

    fill=agecat)) + geom_boxplot()
ggplot(subset(data1, Clicks>0), aes(x=Clicks, colour=agecat))
  + geom_density()

# create categories
data1$scode[data1$Impressions==0] <- "NoImps"
data1$scode[data1$Impressions >0] <- "Imps"
data1$scode[data1$Clicks >0] <- "Clicks"

# Convert the column to a factor
data1$scode <- factor(data1$scode)
head(data1)

#look at levels
cLen <- function(x){c(length(x))}
etable<-summaryBy(Impressions~scode+Gender+agecat,
  data = data1, FUN=cLen)

```

Hint for doing the rest: don't read all the datasets into memory. Once you've perfected your code for one day, read the datasets in one at a time, process them, output any relevant metrics and variables, and store them in a dataframe; then remove the dataset before reading in the next one. This is to get you thinking about how to handle data sharded across multiple machines.

On Coding

In a May 2013 op-ed piece, "How to be a Woman Programmer," Ellen Ullman describes quite well what it takes to be a programmer (setting aside for now the woman part):

"The first requirement for programming is a passion for the work, a deep need to probe the mysterious space between human thoughts and what a machine can understand; between human desires and how machines might satisfy them.

The second requirement is a high tolerance for failure. Programming is the art of algorithm design and the craft of debugging errant code. In the words of the great John Backus, inventor of the Fortran programming language: *You need the willingness to fail all the time. You have to generate many ideas and then you have to work very hard only to discover that they don't work. And you keep doing that over and over until you find one that does work.*"

The Data Science Process

Let's put it all together into what we define as the data science process. The more examples you see of people doing data science, the more you'll find that they fit into the general framework shown in [Figure 2-2](#). As we go through the book, we'll revisit stages of this process and examples of it in different ways.

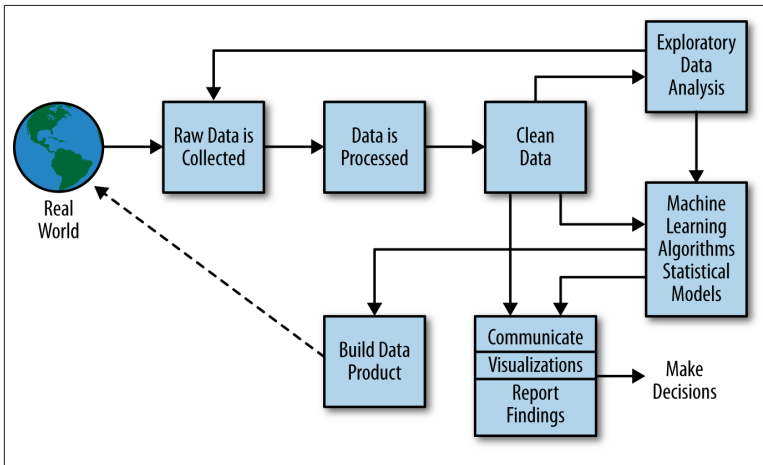


Figure 2-2. The data science process

First we have the Real World. Inside the Real World are lots of people busy at various activities. Some people are using Google+, others are competing in the Olympics; there are spammers sending spam, and there are people getting their blood drawn. Say we have data on one of these things.

Specifically, we'll start with raw data—logs, Olympics records, Enron employee emails, or recorded genetic material (note there are lots of aspects to these activities already lost even when we have that raw data). We want to process this to make it clean for analysis. So we build and use pipelines of data munging: joining, scraping, wrangling, or whatever you want to call it. To do this we use tools such as Python, shell scripts, R, or SQL, or all of the above.

Eventually we get the data down to a nice format, like something with columns:

```
name | event | year | gender | event time
```




This is where you typically *start* in a standard statistics class, with a clean, orderly dataset. But it's not where you typically start in the real world.

Once we have this clean dataset, we should be doing some kind of EDA. In the course of doing EDA, we may realize that it isn't actually clean because of duplicates, missing values, absurd outliers, and data that wasn't actually logged or incorrectly logged. If that's the case, we may have to go back to collect more data, or spend more time cleaning the dataset.

Next, we design our model to use some algorithm like k-nearest neighbor (k-NN), linear regression, Naive Bayes, or something else. The model we choose depends on the type of problem we're trying to solve, of course, which could be a classification problem, a prediction problem, or a basic description problem.

We then can interpret, visualize, report, or communicate our results. This could take the form of reporting the results up to our boss or coworkers, or publishing a paper in a journal and going out and giving academic talks about it.

Alternatively, our goal may be to build or prototype a "data product"; e.g., a spam classifier, or a search ranking algorithm, or a recommendation system. Now the key here that makes data science special and distinct from statistics is that this data product then *gets incorporated back* into the real world, and users interact with that product, and that generates more data, which creates a feedback loop.

This is very different from predicting the weather, say, where your model doesn't influence the outcome at all. For example, you might predict it will rain next week, and unless you have some powers we don't know about, you're not going to *cause* it to rain. But if you instead build a recommendation system that generates evidence that "lots of people love this book," say, then you will know that you caused that feedback loop.

Take this loop into account in any analysis you do by adjusting for any biases your model caused. Your models are not just predicting the future, but *causing* it!

A data product that is productionized and that users interact with is at one extreme and the weather is at the other, but regardless of the

type of data you work with and the “data product” that gets built on top of it—be it public policy determined by a statistical model, health insurance, or election polls that get widely reported and perhaps influence viewer opinions—you should consider the extent to which your model is influencing the very phenomenon that you are trying to observe and understand.

A Data Scientist’s Role in This Process

This model so far seems to suggest this will all magically happen without human intervention. By “human” here, we mean “data scientist.” Someone has to make the decisions about what data to collect, and why. That person needs to be formulating questions and hypotheses and making a plan for how the problem will be attacked. And that someone is the data scientist or our beloved data science team.

Let’s revise or at least add an overlay to make clear that the data scientist needs to be involved in this process throughout, meaning they are involved in the actual coding as well as in the higher-level process, as shown in [Figure 2-3](#).

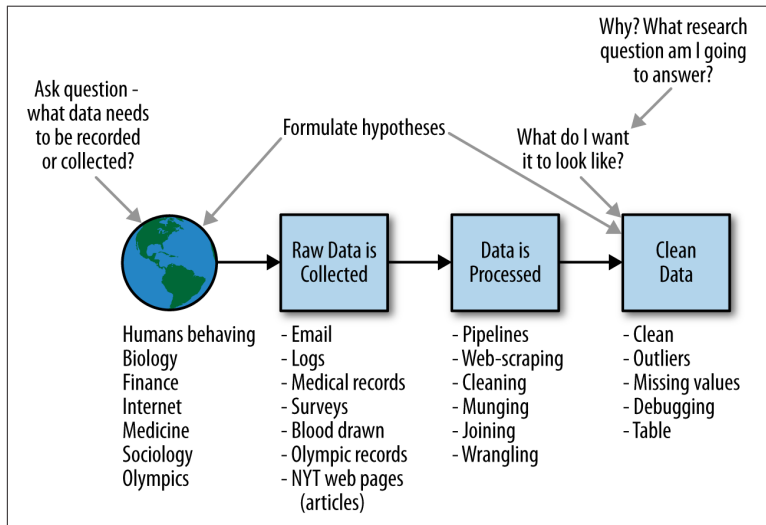


Figure 2-3. The data scientist is involved in every part of this process

Connection to the Scientific Method

We can think of the data science process as an extension of or variation of the scientific method:

- Ask a question.
- Do background research.
- Construct a hypothesis.
- Test your hypothesis by doing an experiment.
- Analyze your data and draw a conclusion.
- Communicate your results.

In both the data science process and the scientific method, not every problem requires one to go through all the steps, but almost all problems can be solved with some combination of the stages. For example, if your end goal is a data visualization (which itself could be thought of as a data product), it's possible you might not do any machine learning or statistical modeling, but you'd want to get all the way to a clean dataset, do some exploratory analysis, and then create the visualization.

Thought Experiment: How Would You Simulate Chaos?

Most data problems start out with a certain amount of dirty data, ill-defined questions, and urgency. As data scientists we are, in a sense, attempting to create order from chaos. The class took a break from the lecture to discuss how they'd simulate chaos. Here are some ideas from the discussion:

- A Lorenzian water wheel, which is a Ferris wheel-type contraption with equally spaced buckets of water that rotate around in a circle. Now imagine water being dripped into the system at the very top. Each bucket has a leak, so some water escapes into whatever bucket is directly below the drip. Depending on the rate of the water coming in, this system exhibits a chaotic process that depends on molecular-level interactions of water molecules on the sides of the buckets. Read more about it [in this associated Wikipedia article](#).

- Many systems can exhibit inherent chaos. Philippe M. Binder and Roderick V. Jensen have written a paper entitled “**Simulating chaotic behavior with finite-state machines**”, which is about digital computer simulations of chaos.
- An interdisciplinary program involving M.I.T., Harvard, and Tufts involved teaching a technique that was entitled “**Simulating chaos to teach order**”. They simulated an emergency on the border between Chad and Sudan’s troubled Darfur region, with students acting as members of Doctors Without Borders, International Medical Corps, and other humanitarian agencies.
- See also Joel Gascoigne’s related essay, “**Creating order from chaos in a startup**”.

Instructor Notes

1. Being a data scientist in an organization is often a chaotic experience, and it’s the data scientist’s job to try to create order from that chaos. So I wanted to simulate that chaotic experience for my students throughout the semester. But I also wanted them to know that things were going to be slightly chaotic for a pedagogical reason, and not due to my ineptitude!
2. I wanted to draw out different interpretations of the word “chaos” as a means to think about the importance of vocabulary, and the difficulties caused in communication when people either don’t know what a word means, or have different ideas of what the word means. Data scientists might be communicating with domain experts who don’t really understand what “logistic regression” means, say, but will pretend to know because they don’t want to appear stupid, or because they think they ought to know, and therefore don’t ask. But then the whole conversation is not really a successful communication if the two people talking don’t really understand what they’re talking about. Similarly, the data scientists ought to be asking questions to make sure they understand the terminology the domain expert is using (be it an astrophysicist, a social networking expert, or a climatologist). There’s nothing wrong with not knowing what a word means, but there is something wrong with not asking! You will likely find that asking clarifying questions about vocabulary gets you even more insight into the underlying data problem.

3. Simulation is a useful technique in data science. It can be useful practice to simulate fake datasets from a model to understand the generative process better, for example, and also to debug code.

Case Study: RealDirect

Doug Perlson, the CEO of **RealDirect**, has a background in real estate law, startups, and online advertising. His goal with RealDirect is to use all the data he can access about real estate to improve the way people sell and buy houses.

Normally, people sell their homes about once every seven years, and they do so with the help of professional brokers and current data. But there's a problem both with the broker system and the data quality. RealDirect addresses both of them.

First, the brokers. They are typically “free agents” operating on their own—think of them as home sales consultants. This means that they guard their data aggressively, and the really good ones have lots of experience. But in the grand scheme of things, that really means they have only slightly more data than the inexperienced brokers.

RealDirect is addressing this problem by hiring a team of licensed real-estate agents who work together and pool their knowledge. To accomplish this, it built an interface for sellers, giving them useful data-driven tips on how to sell their house. It also uses interaction data to give real-time recommendations on what to do next.

The team of brokers also become data experts, learning to use information-collecting tools to keep tabs on new and relevant data or to access publicly available information. For example, you can now get data on co-op (a certain kind of apartment in NYC) sales, but that's a relatively recent change.

One problem with publicly available data is that it's old news—there's a three-month lag between a sale and when the data about that sale is available. RealDirect is working on real-time feeds on things like when people start searching for a home, what the initial offer is, the time between offer and close, and how people search for a home online.

Ultimately, good information helps both the buyer and the seller. At least if they're honest.

How Does RealDirect Make Money?

First, it offers a subscription to sellers—about \$395 a month—to access the selling tools. Second, it allows sellers to use RealDirect’s agents at a reduced commission, typically 2% of the sale instead of the usual 2.5% or 3%. This is where the magic of data pooling comes in: it allows RealDirect to take a smaller commission because it’s more optimized, and therefore gets more volume.

The site itself is best thought of as a platform for buyers and sellers to manage their sale or purchase process. There are statuses for each person on site: active, offer made, offer rejected, showing, in contract, etc. Based on your status, different actions are suggested by the software.

There are some challenges they have to deal with as well, of course. First off, there’s a law in New York that says you can’t show all the current housing listings unless those listings reside behind a registration wall, so RealDirect requires registration. On the one hand, this is an obstacle for buyers, but serious buyers are likely willing to do it. Moreover, places that don’t require registration, like [Zillow](#), aren’t true competitors to RealDirect because they are merely showing listings without providing any additional service. Doug pointed out that you also need to register to use [Pinterest](#), and it has tons of users in spite of this.

RealDirect comprises licensed brokers in various established realtor associations, but even so it has had its share of hate mail from realtors who don’t appreciate its approach to cutting commission costs. In this sense, RealDirect is breaking directly into a guild. On the other hand, if a realtor refused to show houses because they are being sold on RealDirect, the potential buyers would see those listings elsewhere and complain. So the traditional brokers have little choice but to deal with RealDirect even if they don’t like it. In other words, the listings themselves are sufficiently transparent so that the traditional brokers can’t get away with keeping their buyers away from these houses.

Doug talked about key issues that a buyer might care about—nearby parks, subway, and schools, as well as the comparison of prices per square foot of apartments sold in the same building or block. This is the kind of data they want to increasingly cover as part of the service of RealDirect.

Exercise: RealDirect Data Strategy

You have been hired as chief data scientist at *realdirect.com*, and report directly to the CEO. The company (hypothetically) does not yet have its data plan in place. It's looking to you to come up with a data strategy. Here are a couple ways you could begin to approach this problem:

1. Explore its existing website, thinking about how buyers and sellers would navigate through it, and how the website is structured/organized. Try to understand the existing business model, and think about how analysis of RealDirect user-behavior data could be used to inform decision-making and product development. Come up with a list of research questions you think could be answered by data:
 - What data would you advise the engineers log and what would your ideal datasets look like?
 - How would data be used for reporting and monitoring product usage?
 - How would data be built back into the product/website?
2. Because there is no data yet for you to analyze (typical in a start-up when its still building its product), you should get some auxiliary data to help gain intuition about this market. For example, go to https://github.com/oreillymedia/doing_data_science. Click on Rolling Sales Update (after the fifth paragraph).

You can use any or all of the datasets here—start with Manhattan August, 2012–August 2013.

- First challenge: load in and clean up the data. Next, conduct exploratory data analysis in order to find out where there are outliers or missing values, decide how you will treat them, make sure the dates are formatted correctly, make sure values you think are numerical are being treated as such, etc.
 - Once the data is in good shape, conduct exploratory data analysis to visualize and make comparisons (i) across neighborhoods, and (ii) across time. If you have time, start looking for meaningful patterns in this dataset.
3. Summarize your findings in a brief report aimed at the CEO.

4. Being the “data scientist” often involves speaking to people who aren’t also data scientists, so it would be ideal to have a set of communication strategies for getting to the information you need about the data. Can you think of any other people you should talk to?
5. Most of you are not “domain experts” in real estate or online businesses.
 - Does stepping out of your comfort zone and figuring out how you would go about “collecting data” in a different setting give you insight into how you do it in your own field?
 - Sometimes “domain experts” have their own set of vocabulary. Did Doug use vocabulary specific to his domain that you didn’t understand (“comps,” “open houses,” “CPC”)? Sometimes if you don’t understand vocabulary that an expert is using, it can prevent you from understanding the problem. It’s good to get in the habit of asking questions because eventually you will get to something you do understand. This involves persistence and is a habit to cultivate.
6. Doug mentioned the company didn’t necessarily have a data strategy. There is no industry standard for creating one. As you work through this assignment, think about whether there is a set of best practices you would recommend with respect to developing a data strategy for an online business, or in your own domain.

Sample R code

Here’s some sample R code that takes the Brooklyn housing data in the preceding exercise, and cleans and explores it a bit. (The exercise asks you to do this for Manhattan.)

```
# Author: Benjamin Reddy

require(gdata)
bk <- read.xls("rollingsales_brooklyn.xls", pattern="BOROUGH")
head(bk)
summary(bk)

bk$SALE.PRICE.N <- as.numeric(gsub("[^[:digit:]]", "",
                                bk$SALE.PRICE))
count(is.na(bk$SALE.PRICE.N))

names(bk) <- tolower(names(bk))
```



```

## clean/format the data with regular expressions
bk$gross.sqft <- as.numeric(gsub("[^[:digit:]]", "",
                               bk$gross.square.feet))
bk$land.sqft <- as.numeric(gsub("[^[:digit:]]", "",
                               bk$land.square.feet))

bk$sale.date <- as.Date(bk$sale.date)
bk$year.built <- as.numeric(as.character(bk$year.built))

## do a bit of exploration to make sure there's not anything
## weird going on with sale prices
attach(bk)

hist(sale.price.n)
hist(sale.price.n[sale.price.n>0])
hist(gross.sqft[sale.price.n==0])

detach(bk)

## keep only the actual sales
bk.sale <- bk[bk$sale.price.n!=0,]

plot(bk.sale$gross.sqft,bk.sale$sale.price.n)
plot(log(bk.sale$gross.sqft),log(bk.sale$sale.price.n))

## for now, let's look at 1-, 2-, and 3-family homes
bk.homes <- bk.sale[which(grepl("FAMILY",
                               bk.sale$building.class.category)),]
plot(log(bk.homes$gross.sqft),log(bk.homes$sale.price.n))

bk.homes[which(bk.homes$sale.price.n<100000),]
  [order(bk.homes[which(bk.homes$sale.price.n<100000),]
        $sale.price.n),]

## remove outliers that seem like they weren't actual sales
bk.homes$outliers <- (log(bk.homes$sale.price.n) <=5) + 0
bk.homes <- bk.homes[which(bk.homes$outliers==0),]

plot(log(bk.homes$gross.sqft),log(bk.homes$sale.price.n))

```

Algorithms

In the previous chapter we discussed in general how models are used in data science. In this chapter, we're going to be diving into algorithms.

An algorithm is a procedure or set of steps or rules to accomplish a task. Algorithms are one of the fundamental concepts in, or building blocks of, computer science: the basis of the design of elegant and efficient code, data preparation and processing, and software engineering.

Some of the basic types of tasks that algorithms can solve are sorting, searching, and graph-based computational problems. Although a given task such as sorting a list of objects could be handled by multiple possible algorithms, there is some notion of “best” as measured by efficiency and computational time, which matters especially when you're dealing with massive amounts of data and building consumer-facing products.

Efficient algorithms that work sequentially or in parallel are the basis of pipelines to process and prepare data. With respect to data science, there are at least three classes of algorithms one should be aware of:

1. Data munging, preparation, and processing algorithms, such as sorting, MapReduce, or Pregel.

We would characterize these types of algorithms as data engineering, and while we devote a chapter to this, it's not the emphasis of this book. This is not to say that you won't be doing data wrangling and munging—just that we don't emphasize the algorithmic aspect of it.

2. Optimization algorithms for parameter estimation, including Stochastic Gradient Descent, Newton’s Method, and Least Squares. We mention these types of algorithms throughout the book, and they underlie many R functions.
3. Machine learning algorithms are a large part of this book, and we discuss these more next.

Machine Learning Algorithms

Machine learning algorithms are largely used to predict, classify, or cluster.

Wait! Back in the previous chapter, didn’t we already say modeling could be used to predict or classify? Yes. Here’s where some lines have been drawn that can make things a bit confusing, and it’s worth understanding who drew those lines.

Statistical *modeling* came out of statistics departments, and machine learning *algorithms* came out of computer science departments. Certain methods and techniques are considered to be part of both, and you’ll see that we often use the words somewhat interchangeably.

You’ll find some of the methods in this book, such as linear regression, in machine learning books as well as intro to statistics books. It’s not necessarily useful to argue over who the rightful owner is of these methods, but it’s worth pointing out here that it can get a little vague or ambiguous about what the actual difference is.

In general, machine learning algorithms that are the basis of artificial intelligence (AI) such as image recognition, speech recognition, recommendation systems, ranking and personalization of content—often the basis of data products—are not usually part of a core statistics curriculum or department. They aren’t generally designed to infer the underlying *generative process* (e.g., to model something), but rather to predict or classify with the most accuracy.

These differences in methods reflect in cultural differences in the approaches of machine learners and statisticians that Rachel observed at Google, and at industry conferences. Of course, data scientists can and should use both approaches.

There are some broad generalizations to consider:

Interpreting parameters

Statisticians think of the parameters in their linear regression models as having real-world interpretations, and typically want to be able to find meaning in behavior or describe the real-world phenomenon corresponding to those parameters. Whereas a software engineer or computer scientist might be wanting to build their linear regression algorithm into production-level code, and the predictive model is what is known as a *black box algorithm*, they don't generally focus on the interpretation of the parameters. If they do, it is with the goal of handtuning them in order to optimize *predictive power*.

Confidence intervals

Statisticians provide confidence intervals and posterior distributions for parameters and estimators, and are interested in capturing the variability or uncertainty of the parameters. Many machine learning algorithms, such as k-means or k-nearest neighbors (which we cover a bit later in this chapter), don't have a notion of confidence intervals or uncertainty.

The role of explicit assumptions

Statistical models make explicit assumptions about data-generating processes and distributions, and you use the data to estimate parameters. Nonparametric solutions, like we'll see later in this chapter, don't make any assumptions about probability distributions, or they are implicit.

We say the following lovingly and with respect: statisticians have chosen to spend their lives investigating uncertainty, and they're never 100% confident about anything. Software engineers like to build things. They want to build models that predict the best they can, but there are no concerns about uncertainty—just build it! At companies like Facebook or Google, the philosophy is to build and iterate often. If something breaks, it can be fixed. A data scientist who somehow manages to find a balance between the statistical and computer science approaches, and to find value in both these ways of being, can thrive. Data scientists are the multicultural statistician-computer scientist hybrid, so we're not tied to any one way of thinking over another; they both have value. We'll sum up our take on this with guest speaker Josh Wills' (Chapter 13) well-tweeted quote:

Data scientist (noun): Person who is better at statistics than any software engineer and better at software engineering than any statistician.

— Josh Wills

Three Basic Algorithms

Many business or real-world problems that can be solved with data can be thought of as *classification* and *prediction* problems when we express them mathematically. Happily, a whole host of models and algorithms can be used to classify and predict.

Your real challenge as a data scientist, once you've become familiar with how to implement them, is understanding which ones to use depending on the context of the problem and the underlying assumptions. This partially comes with experience—you start seeing enough problems that you start thinking, “Ah, this is a classification problem with a binary outcome” or, “This is a classification problem, but oddly I don't even have any labels” and you know what to do. (In the first case, you could use logistic regression or Naive Bayes, and in the second you could start with k-means—more on all these shortly!)

Initially, though, when you hear about these methods in isolation, it takes some effort on your part as a student or learner to think, “In the real world, how do I know that this algorithm is the solution to the problem I'm trying to solve?”

It's a real mistake to be the type of person walking around with a hammer looking for a nail to bang: “I know linear regression, so I'm going to try to solve every problem I encounter with it.” Don't do that. Instead, try to understand the context of the problem, and the attributes it has *as a problem*. Think of those in mathematical terms, and then think about the algorithms you know and how they map to this type of problem.

If you're not sure, it's good to talk it through with someone who does. So ask a coworker, head to a meetup group, or start one in your area! Also, maintain the attitude that it's *not obvious* what to do and that's what makes it a problem, and so you're going to approach it circumspectly and methodically. You don't have to be the know-it-all in the room who says, “Well, *obviously* we should use linear regression with a penalty function for regularization,” even if that seems to you the right approach.

We're saying all this because one of the unfortunate aspects of textbooks is they often give you a bunch of techniques and then problems that tell you *which* method to use that solves the problem (e.g., use linear regression to predict height from weight). Yes, implementing and understanding linear regression the first few times is not obvious, so you need practice with that, but it needs to be addressed that the real challenge once you have mastery over technique is *knowing when to use linear regression in the first place*.

We're not going to give a comprehensive overview of *all* possible machine learning algorithms, because that would make this a machine learning book, and there are already plenty of those.

Having said that, in this chapter we'll introduce three basic algorithms now and introduce others throughout the book in context. By the end of the book, you should feel more confident about your ability to learn new algorithms so that you can pick them up along the way as problems require them.

We'll also do our best to demonstrate the thought processes of data scientists who had to figure out which algorithm to use in context and why, but it's also upon you as a student and learner to *force yourself* to think about what the attributes of the problem were that made a given algorithm the right algorithm to use.

With that said, we still need to give you some basic tools to use, so we'll start with linear regression, k-nearest neighbors (k-NN), and k-means. In addition to what was just said about trying to understand the attributes of problems that could use these as solutions, look at these three algorithms from the perspective of: what patterns can we as humans see in the data with our eyes that we'd like to be able to automate with a machine, especially taking into account that as the data gets more complex, we can't see these patterns?

Linear Regression

One of the most common statistical methods is linear regression. At its most basic, it's used when you want to express the mathematical relationship between two variables or attributes. When you use it, you are making the assumption that there is a *linear* relationship between an outcome variable (sometimes also called the response variable, dependent variable, or label) and a predictor (sometimes also called an independent variable, explanatory variable, or feature); or between

one variable and several other variables, in which case you're *modeling* the relationship as having a linear structure.

WTF. So Is It an Algorithm or a Model?

While we tried to make a distinction between the two earlier, we admit the colloquial use of the words “model” and “algorithm” gets confusing because the two words seem to be used interchangeably when their actual definitions are not the same thing at all. In the purest sense, an algorithm is a set of rules or steps to follow to accomplish some task, and a model is an attempt to describe or capture the world. These two seem obviously different, so it seems the distinction should be obvious. Unfortunately, it isn't. For example, regression can be described as a statistical model as well as a machine learning algorithm. You'll waste your time trying to get people to discuss this with any precision.

In some ways this is a historical artifact of statistics and computer science communities developing methods and techniques in parallel and using different words for the same methods. The consequence of this is that the distinction between machine learning and statistical modeling is muddy. Some methods (for example, k-means, discussed in the next section) we might call an *algorithm* because it's a series of computational steps used to cluster or classify objects—on the other hand, k-means can be reinterpreted as a special case of a Gaussian mixture *model*. The net result is that colloquially, people use the terms algorithm and model interchangeably when it comes to a lot of these methods, so try not to let it worry you. (Though it bothers us, too.)

Assuming that there is a *linear* relationship between an outcome variable and a predictor is a big assumption, but it's also the simplest one you *can* make—linear functions are more basic than nonlinear ones in a mathematical sense—so in that sense it's a good starting point.

In some cases, it makes sense that changes in one variable correlate linearly with changes in another variable. For example, it makes sense that the more umbrellas you sell, the more money you make. In those cases you can feel good about the linearity assumption. Other times, it's harder to justify the assumption of linearity except locally: in the spirit of calculus, everything can be approximated by line segments as long as functions are continuous.

Let's back up. Why would you even want to build a linear model in the first place? You might want to use this relationship to *predict* future outcomes, or you might want to understand or *describe* the relationship to get a grasp on the situation. Let's say you're studying the relationship between a company's sales and how much that company spends on advertising, or the number of friends someone has on a social networking site and the time that person spends on that site daily. These are all numerical outcomes, which mean linear regression would be a wise choice, at least for a first pass at your problem.

One entry point for thinking about linear regression is to think about deterministic lines first. We learned back in grade school that we could describe a line with a slope and an intercept, $y = f(x) = \beta_0 + \beta_1 * x$. But the setting there was always deterministic.

Even for the most mathematically sophisticated among us, if you haven't done it before, it's a new mindset to start thinking about stochastic functions. We still have the same components: points listed out explicitly in a table (or as tuples), and functions represented in equation form or plotted on a graph. So let's build up to linear regression starting from a deterministic function.

Example 1. Overly simplistic example to start. Suppose you run a social networking site that charges a monthly subscription fee of \$25, and that this is your only source of revenue. Each month you collect data and count your number of users and total revenue. You've done this daily over the course of two years, recording it all in a spreadsheet. You could express this data as a series of points. Here are the first four:

$$S = \{(x, y) = (1, 25), (10, 250), (100, 2500), (200, 5000)\}$$

If you showed this to someone else who didn't even know how much you charged or anything about your business model (what kind of friend wasn't paying attention to your business model?!), they might notice that there's a clear relationship enjoyed by all of these points, namely $y = 25x$. They likely could do this in their head, in which case they figured out that:

- There's a linear pattern.
- The coefficient relating x and y is 25.
- It seems deterministic.

You can even plot it as in [Figure 3-1](#) to verify they were right (even though you knew they were because you made the business model in the first place). It's a line!

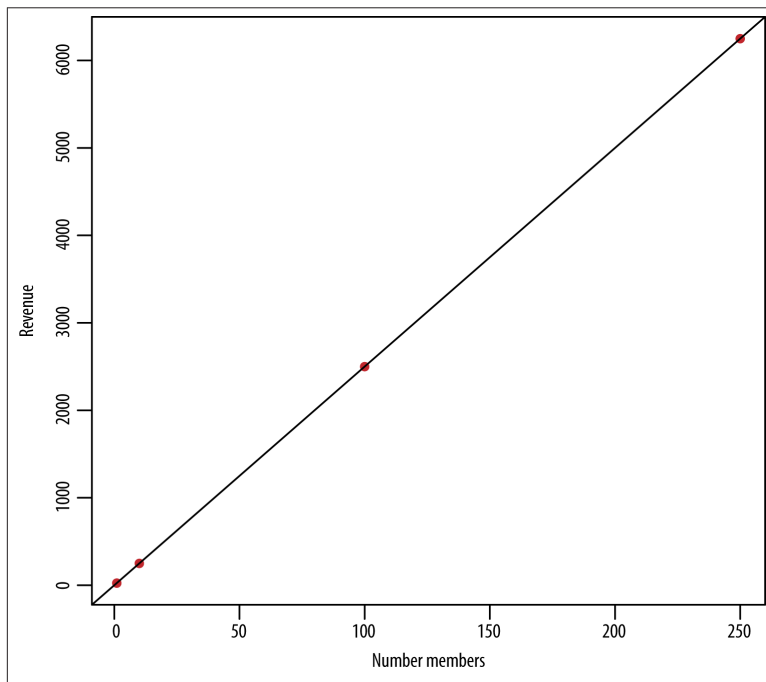


Figure 3-1. An obvious linear pattern

Example 2. Looking at data at the user level. Say you have a dataset *keyed* by user (meaning each row contains data for a single user), and the columns represent user behavior on a social networking site over a period of a week. Let's say you feel comfortable that the data is clean at this stage and that you have on the order of hundreds of thousands of users. The names of the columns are `total_num_friends`, `total_new_friends_this_week`, `num_visits`, `time_spent`, `number_apps_downloaded`, `number_ads_shown`, `gender`, `age`, and so on. During the course of your exploratory data analysis, you've randomly sampled 100 users to keep it simple, and you plot pairs of these variables, for example, $x = \text{total_new_friends}$ and $y = \text{time_spent}$ (in seconds). The business context might be that eventually you want to be able to promise advertisers who bid for space on your website in advance a certain number of users, so you want to be able to forecast

number of users several days or weeks in advance. But for now, you are simply trying to build intuition and understand your dataset.

You eyeball the first few rows and see:

```
7 276
3 43
4 82
6 136
10 417
9 269
```

Now, your brain can't figure out what's going on by just looking at them (and your friend's brain probably can't, either). They're in no obvious particular order, and there are a lot of them. So you try to plot it as in [Figure 3-2](#).

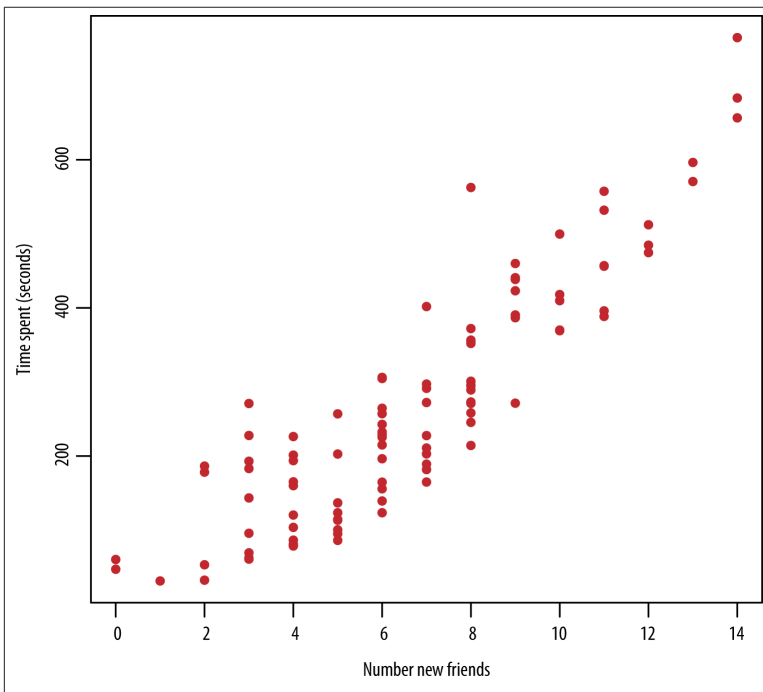


Figure 3-2. Looking kind of linear

It looks like there's *kind of* a linear relationship here, and it makes sense; the more new friends you have, the more time you might spend on the site. But how can you figure out how to describe that relationship? Let's also point out that there is no perfectly *deterministic* relationship between number of new friends and time spent on the site, but it makes sense that there is an *association* between these two variables.

Start by writing something down

There are two things you want to capture in the model. The first is the *trend* and the second is the *variation*. We'll start first with the trend.

First, let's start by assuming there actually *is* a relationship and that it's linear. It's the best you can do at this point.

There are many lines that look more or less like they might work, as shown in [Figure 3-3](#).

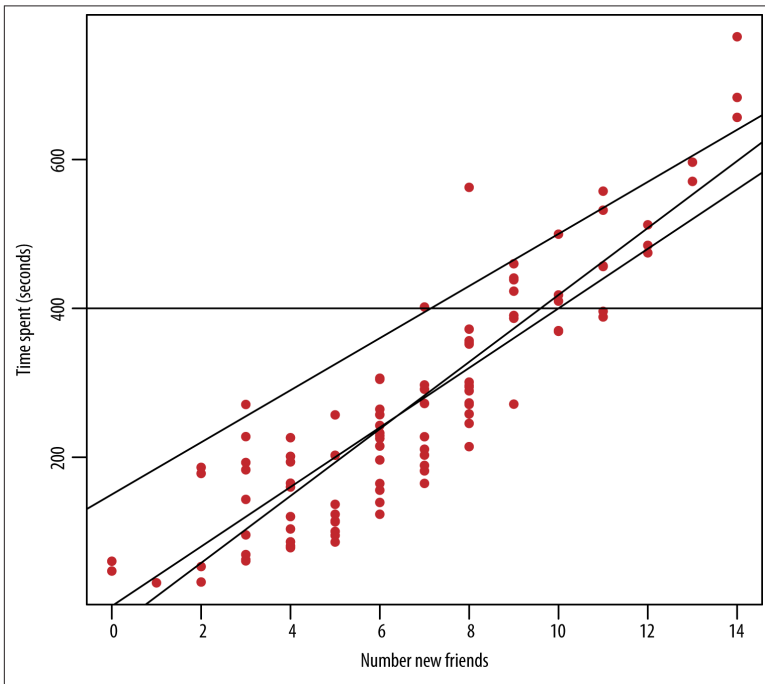


Figure 3-3. Which line is the best fit?

So how do you pick which one?

Because you're assuming a linear relationship, start your model by assuming the functional form to be:

$$y = \beta_0 + \beta_1 x$$

Now your job is to find the best choices for β_0 and β_1 using the observed data to estimate them: $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$.

Writing this with matrix notation results in this:

$$y = x \cdot \beta$$

There you go: you've written down your model. Now the rest is *fitting* the model.

Fitting the model

So, how do you calculate β ? The intuition behind linear regression is that you want to find the line that minimizes the distance between all the points and the line.

Many lines look approximately correct, but your goal is to find the optimal one. Optimal could mean different things, but let's start with optimal to mean the line that, on average, is closest to all the points. But what does *closest* mean here?

Look at [Figure 3-4](#). Linear regression seeks to find the line that minimizes the sum of the squares of the vertical distances between the approximated or predicted \hat{y}_i s and the observed y_i s. You do this because you want to minimize your prediction errors. This method is called *least squares* estimation.

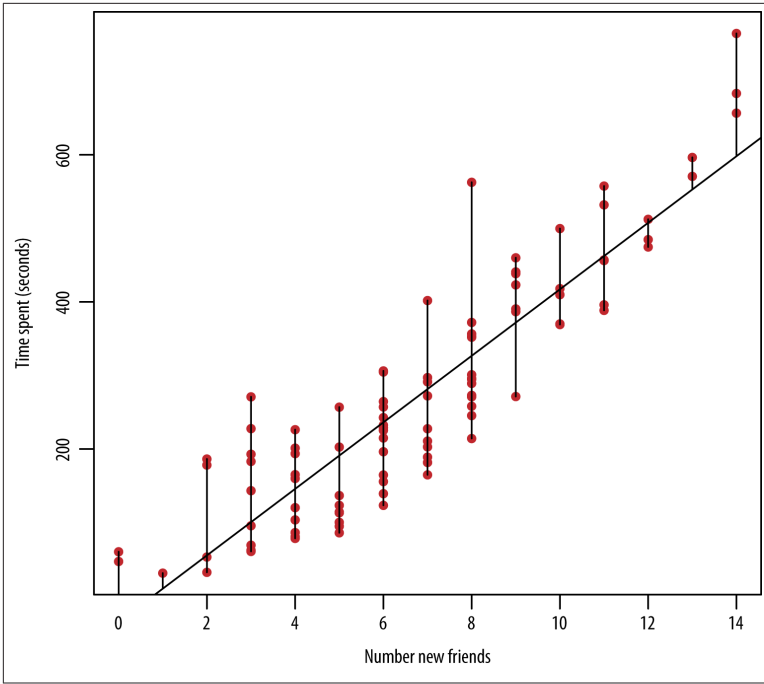


Figure 3-4. The line closest to all the points

To find this line, you’ll define the “residual sum of squares” (RSS), denoted $RSS(\beta)$, to be:

$$RSS(\beta) = \sum_i (y_i - \beta x_i)^2$$

where i ranges over the various data points. It is the sum of all the squared vertical distances between the observed points and any given line. Note this is a function of β and you want to optimize with respect to β to find the optimal line.

To minimize $RSS(\beta) = (y - \beta x)^t (y - \beta x)$, differentiate it with respect to β and set it equal to zero, then solve for β . This results in:

$$\hat{\beta} = (x^t x)^{-1} x^t y$$

Here the little “hat” symbol on top of the β is there to indicate that it’s the *estimator* for β . You don’t know the true value of β ; all you have is the observed data, which you plug into the estimator to get an estimate.

To actually fit this, to get the β s, all you need is one line of R code where you’ve got a column of y’s and a (single) column of x’s:

```
model <- lm(y ~ x)
```

So for the example where the first few rows of the data were:

```
x y
7 276
3 43
4 82
6 136
10 417
9 269
```

The R code for this would be:

```
> model <- lm (y~x)
> model

Call:
lm(formula = y ~ x)

Coefficients:
(Intercept)          x
      -32.08         45.92

> coefs <- coef(model)
> plot(x, y, pch=20,col="red", xlab="Number new friends",
      ylab="Time spent (seconds)")
> abline(coefs[1],coefs[2])
```

And the estimated line is $\hat{y} = -32.08 + 45.92x$, which you’re welcome to round to $\hat{y} = -32 + 46x$, and the corresponding plot looks like the lefthand side of [Figure 3-5](#).

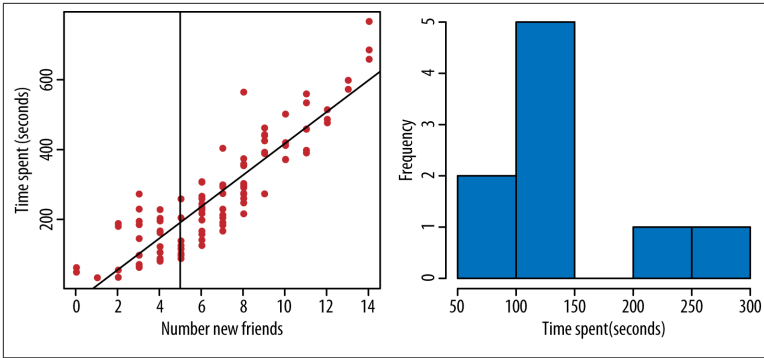


Figure 3-5. On the left is the fitted line. We can see that for any fixed value, say 5, the values for y vary. For people with 5 new friends, we display their time spent in the plot on the right.

But it's up to you, the data scientist, whether you think you'd actually want to use this linear model to describe the relationship or predict new outcomes. If a new x -value of 5 came in, meaning the user had five new friends, how confident are you in the output value of $-32.08 + 45.92 \times 5 = 195.7$ seconds?

In order to get at this question of confidence, you need to extend your model. You know there's variation among time spent on the site by people with five new friends, meaning you certainly wouldn't make the claim that everyone with five new friends is guaranteed to spend 195.7 seconds on the site. So while you've so far modeled the *trend*, you haven't yet modeled the *variation*.

Extending beyond least squares

Now that you have a *simple linear regression model* down (one output, one predictor) using least squares estimation to estimate your β s, you can build upon that model in three primary ways, described in the upcoming sections:

1. Adding in modeling assumptions about the errors
2. Adding in more predictors
3. Transforming the predictors

Adding in modeling assumptions about the errors. If you use your model to predict y for a given value of x , your prediction is deterministic and doesn't capture the variability in the observed data. See on the

righthand side of Figure 3-5 that for a fixed value of $x=5$, there is variability among the time spent on the site. You want to capture this variability in your model, so you extend your model to:

$$y = \beta_0 + \beta_1 x + \epsilon$$

where the new term ϵ is referred to as *noise*, which is the stuff that you haven't accounted for by the relationships you've figured out so far. It's also called the *error term*— ϵ represents the *actual error*, the difference between the observations and the *true* regression line, which you'll never know and can only estimate with your $\hat{\beta}$ s.

One often makes the modeling assumption that the noise is normally distributed, which is denoted:

$$\epsilon \sim N(0, \sigma^2)$$



Note this is sometimes not a reasonable assumption. If you are dealing with a known fat-tailed distribution, and if your linear model is picking up only a small part of the value of the variable y , then the error terms are likely also fat-tailed. This is the most common situation in financial modeling.

That's not to say we don't use linear regression in finance, though. We just don't attach the "noise is normal" assumption to it.

With the preceding assumption on the distribution of noise, this model is saying that, for any given value of x , the conditional distribution of y given x is $p(y|x) \sim N(\beta_0 + \beta_1 x, \sigma^2)$.

So, for example, among the set of people who had five new friends this week, the amount of the time they spent on the website had a *normal distribution* with a mean of $\beta_0 + \beta_1 * 5$ and a variance of σ^2 , and you're going to estimate your parameters β_0, β_1, σ from the data.

How do you fit this model? How do you get the parameters β_0, β_1, σ from the data?



Turns out that no matter how the ϵ s are distributed, the least squares estimates that you already derived are the optimal estimators for β s because they have the property of being unbiased and of being the minimum variance estimators. If you want to know more about these properties and see a proof for this, we refer you to any good book on statistical inference (for example, *Statistical Inference* by Casella and Berger).

So what can you do with your observed data to estimate the variance of the errors? Now that you have the estimated line, you can see how far away the observed data points are from the line itself, and you can treat these differences, also known as *observed errors* or *residuals*, as observations themselves, or estimates of the actual errors, the ϵ s. Define $e_i = y_i - \hat{y}_i = y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i)$ for $i = 1, \dots, n$.

Then you estimate the variance (σ^2) of ϵ , as:

$$\frac{\sum_i e_i^2}{n-2}$$



Why are we dividing by $n-2$? A natural question. Dividing by $n-2$, rather than just n , produces an *unbiased estimator*. The 2 corresponds to the number of model parameters. Here again, Casella and Berger's book is an excellent resource for more background information.

This is called the *mean squared error* and captures how much the predicted value varies from the observed. *Mean squared error* is a useful quantity for any prediction problem. In regression in particular, it's *also* an estimator for your variance, but it can't always be used or interpreted that way. It appears in the evaluation metrics in the following section.

Evaluation metrics

We asked earlier how confident you would be in these estimates and in your model. You have a couple values in the output of the R function that help you get at the issue of how confident you can be in the estimates: p-values and R-squared. Going back to our model in R, if we

type in `summary(model)`, which is the name we gave to this model, the output would be:

```
summary (model)
Call:
lm(formula = y ~ x)

Residuals:
    Min       1Q   Median       3Q      Max
-121.17  -52.63   -9.72   41.54  356.27

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  -32.083     16.623   -1.93  0.0565 .
x             45.918     2.141   21.45 <2e-16 ***
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 77.47 on 98 degrees of freedom
Multiple R-squared:  0.8244,    Adjusted R-squared:  0.8226
F-statistic: 460 on 1 and 98 DF,  p-value: < 2.2e-16
```

R-squared

$R^2 = 1 - \frac{\sum_i (y_i - \hat{y}_i)^2}{\sum_i (y_i - \bar{y})^2}$. This can be interpreted as the proportion of variance explained by our model. Note that mean squared error is in there getting divided by total error, which is the proportion of variance *unexplained* by our model and we calculate 1 minus that.

p-values

Looking at the output, the estimated β s are in the column marked Estimate. To see the p-values, look at $Pr(>|t|)$. We can interpret the values in this column as follows: We are making a null hypothesis that the β s are zero. For any given β , the p-value captures the probability of observing the data that we observed, and obtaining the test-statistic that we obtained *under the null hypothesis*. This means that if we have a low p-value, it is highly unlikely to observe such a test-statistic under the null hypothesis, and the coefficient is highly likely to be nonzero and therefore significant.

Cross-validation

Another approach to evaluating the model is as follows. Divide our data up into a training set and a test set: 80% in the training and 20% in the test. Fit the model on the training set, then look at the *mean squared error* on the test set and compare it to that on the training set. Make this comparison across sample size as well.

If the mean squared errors are approximately the same, then our model generalizes well and we're not in danger of overfitting. See [Figure 3-6](#) to see what this might look like. This approach is highly recommended.

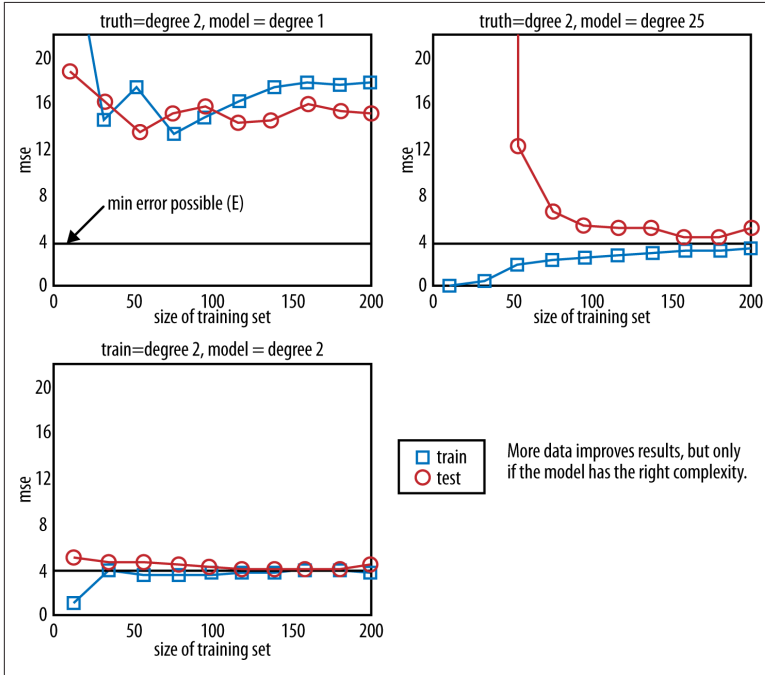


Figure 3-6. Comparing mean squared error in training and test set, taken from a slide of Professor Nando de Freitas; here, the ground truth is known because it came from a dataset with data simulated from a known distribution

Other models for error terms

The mean squared error is an example of what is called a *loss function*. This is the standard one to use in linear regression because it gives us a pretty nice measure of closeness of fit. It has the additional desirable property that by assuming that ϵ s are normally distributed, we can rely on the maximum likelihood principle. There are other loss functions such as one that relies on absolute value rather than squaring. It's also possible to build custom loss functions specific to your particular problem or context, but for now, you're safe with using mean square error.

Adding other predictors. What we just looked at was simple linear regression—one outcome or dependent variable and one predictor. But we can extend this model by building in other predictors, which is called *multiple linear regression*:

$$y = \beta_0 + \beta_1x_1 + \beta_2x_2 + \beta_3x_3 + \epsilon.$$

All the math that we did before holds because we had expressed it in matrix notation, so it was already generalized to give the appropriate estimators for the β . In the example we gave of predicting time spent on the website, the other predictors could be the user's age and gender, for example. We'll explore feature selection more in [Chapter 7](#), which means figuring out which additional predictors you'd want to put in your model. The R code will just be:

```
model <- lm(y ~ x_1 + x_2 + x_3)
```

Or to add in interactions between variables:

```
model <- lm(y ~ x_1 + x_2 + x_3 + x2_*x_3)
```

One key here is to make scatterplots of y against each of the predictors as well as between the predictors, and histograms of $y|x$ for various values of each of the predictors to help build intuition. As with simple linear regression, you can use the same methods to evaluate your model as described earlier: looking at R^2 , p-values, and using training and testing sets.

Transformations. Going back to one x predicting one y , why did we assume a linear relationship? Instead, maybe, a better model would be a polynomial relationship like this:

$$y = \beta_0 + \beta_1x + \beta_2x^2 + \beta_3x^3$$

Wait, but isn't this *linear* regression? Last time we checked, polynomials weren't linear. To think of it as *linear*, you transform or create new variables—for example, $z = x^2$ —and build a regression model based on z . Other common transformations are to take the log or to pick a threshold and turn it into a binary predictor instead.

If you look at the plot of time spent versus number friends, the shape looks a little bit curvy. You could potentially explore this further by building up a model and checking to see whether this yields an improvement.

What you're facing here, though, is one of the biggest challenges for a modeler: you never know the truth. It's possible that the true model is quadratic, but you're assuming linearity or vice versa. You do your best to evaluate the model as discussed earlier, but you'll never *really* know if you're right. More and more data can sometimes help in this regard as well.

Review

Let's review the assumptions we made when we built and fit our model:

- Linearity
- Error terms normally distributed with mean 0
- Error terms independent of each other
- Error terms have constant variance across values of x
- The predictors we're using are the *right* predictors

When and why do we perform linear regression? Mostly for two reasons:

- If we want to predict one variable knowing others
- If we want to explain or understand the relationship between two or more things

Exercise

To help understand and explore new concepts, you can simulate fake datasets in R. The advantage of this is that you "play God" because you actually know the underlying truth, and you get to see how good your model is at recovering the truth.

Once you've better understood what's going on with your fake dataset, you can then transfer your understanding to a real one. We'll show you how to simulate a fake dataset here, then we'll give you some ideas for how to explore it further:

```
# Simulating fake data
x_1 <- rnorm(1000,5,7) # from a normal distribution simulate
                        # 1000 values with a mean of 5 and
                        # standard deviation of 7
hist(x_1, col="grey") # plot  $p(x)$ 
true_error <- rnorm(1000,0,2)
true_beta_0 <- 1.1
true_beta_1 <- -8.2
```

```
y <- true_beta_0 + true_beta_1*x_1 + true_error
hist(y) # plot p(y)
plot(x_1,y, pch=20,col="red") # plot p(x,y)
```

1. Build a regression model and see that it recovers the true values of the β s.
2. Simulate another fake variable x_2 that has a Gamma distribution with parameters you pick. Now make the truth be that y is a linear combination of both x_1 and x_2 . Fit a model that only depends on x_1 . Fit a model that only depends on x_2 . Fit a model that uses both. Vary the sample size and make a plot of mean square error of the training set and of the test set versus sample size.
3. Create a new variable, z , that is equal to x_1^2 . Include this as one of the predictors in your model. See what happens when you fit a model that depends on x_1 only and then also on z . Vary the sample size and make a plot of mean square error of the training set and of the test set versus sample size.
4. Play around more by (a) changing parameter values (the true β s), (b) changing the distribution of the true error, and (c) including more predictors in the model with other kinds of probability distributions. (`rnorm()` means randomly generate values from a normal distribution. `rbinom()` does the same for binomial. So look up these functions online and try to find more.)
5. Create scatterplots of all pairs of variables and histograms of single variables.

k-Nearest Neighbors (k-NN)

K-NN is an algorithm that can be used when you have a bunch of objects that have been classified or labeled in some way, and other similar objects that haven't gotten classified or labeled yet, and you want a way to automatically label them.

The objects could be data scientists who have been classified as “sexy” or “not sexy”; or people who have been labeled as “high credit” or “low credit”; or restaurants that have been labeled “five star,” “four star,” “three star,” “two star,” “one star,” or if they really suck, “zero stars.” More seriously, it could be patients who have been classified as “high cancer risk” or “low cancer risk.”

Take a second and think whether or not linear regression would work to solve problems of this type.

OK, so the answer is: it depends. When you use linear regression, the output is a continuous variable. Here the output of your algorithm is going to be a categorical label, so linear regression wouldn't solve the problem as it's described.

However, it's not impossible to solve it with linear regression plus the concept of a "threshold." For example, if you're trying to predict people's credit scores from their ages and incomes, and then picked a threshold such as 700 such that if your prediction for a given person whose age and income you observed was above 700, you'd label their predicted credit as "high," or toss them into a bin labeled "high." Otherwise, you'd throw them into the bin labeled "low." With more thresholds, you could also have more fine-grained categories like "very low," "low," "medium," "high," and "very high."

In order to do it this way, with linear regression you'd have establish the bins as ranges of a continuous outcome. But not everything is on a continuous scale like a credit score. For example, what if your labels are "likely Democrat," "likely Republican," and "likely independent"? What do you do now?

The intuition behind k-NN is to consider the *most similar* other items defined in terms of their attributes, look at their labels, and give the unassigned item the majority vote. If there's a tie, you randomly select among the labels that have tied for first.

So, for example, if you had a bunch of movies that were labeled "thumbs up" or "thumbs down," and you had a movie called "Data Gone Wild" that hadn't been rated yet—you could look at its attributes: length of movie, genre, number of sex scenes, number of Oscar-winning actors in it, and budget. You could then find other movies with similar attributes, look at *their* ratings, and then give "Data Gone Wild" a rating without ever having to watch it.

To automate it, two decisions must be made: first, how do you define *similarity* or closeness? Once you define it, for a given unrated item, you can say how similar *all* the labeled items are to it, and you can take the *most similar* items and call them *neighbors*, who each have a "vote."

This brings you to the second decision: how many neighbors should you look at or "let vote"? This value is *k*, which ultimately you'll choose as the data scientist, and we'll tell you how.

Make sense? Let's try it out with a more realistic example.

Example with credit scores

Say you have the age, income, and a credit category of high or low for a bunch of people and you want to use the age and income to predict the credit label of "high" or "low" for a new person.

For example, here are the first few rows of a dataset, with income represented in thousands:

```
age income credit
69 3 low
66 57 low
49 79 low
49 17 low
58 26 high
44 71 high
```

You can plot people as points on the plane and label people with an empty circle if they have low credit ratings, as shown in [Figure 3-7](#).

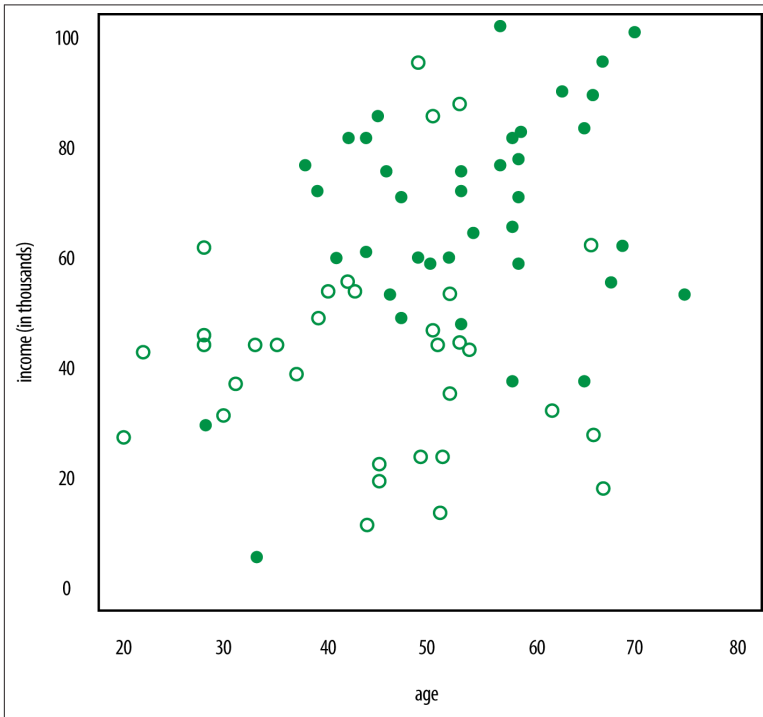


Figure 3-7. Credit rating as a function of age and income

What if a new guy comes in who is 57 years old and who makes \$37,000? What's his likely credit rating label? Look at **Figure 3-8**. Based on the other people near him, what credit score label do you think he should be given? Let's use k-NN to do it automatically.

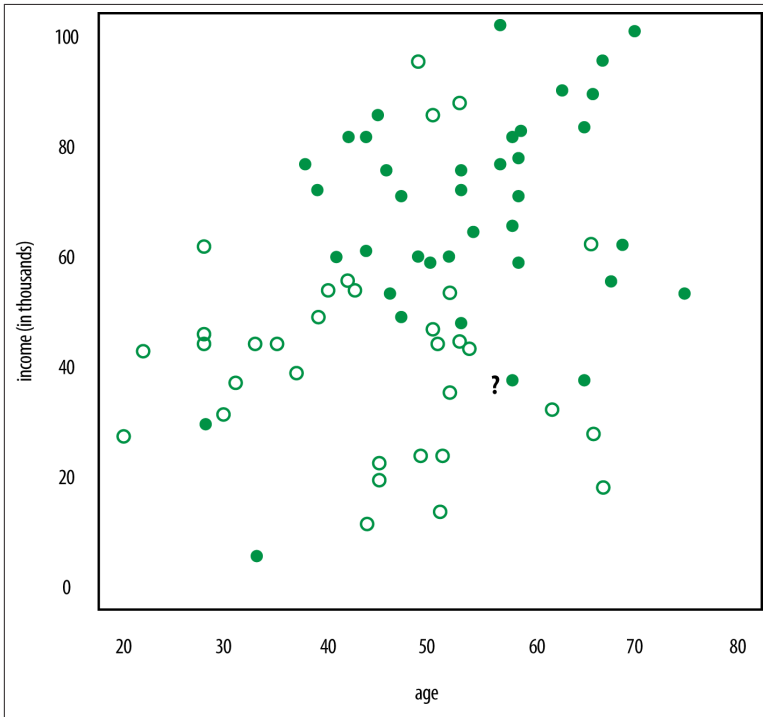


Figure 3-8. What about that guy?

Here's an overview of the process:

1. Decide on your *similarity* or *distance* metric.
2. Split the original labeled dataset into training and test data.
3. Pick an evaluation metric. (Misclassification rate is a good one. We'll explain this more in a bit.)
4. Run k-NN a few times, changing k and checking the evaluation measure.
5. Optimize k by picking the one with the best evaluation measure.
6. Once you've chosen k , use the same training set and now create a new test set with the people's ages and incomes that you have *no*

labels for, and want to predict. In this case, your new test set only has one lonely row, for the 57-year-old.

Similarity or distance metrics

Definitions of “closeness” and similarity vary depending on the context: closeness in social networks could be defined as the number of overlapping friends, for example.

For the sake of our problem of what a neighbor is, we can use Euclidean distance on the plane if the variables are on the same scale. And that can sometimes be a big IF.

Caution: Modeling Danger Ahead!

The scalings question is a really big deal, and if you do it wrong, your model could just suck.

Let’s consider an example: Say you measure age in years, income in dollars, and credit rating as credit scores normally are given—something like SAT scores. Then two people would be represented by triplets such as (25,54000,700) and (35,76000,730). In particular, their “distance” would be completely dominated by the difference in their salaries.

On the other hand, if you instead measured salary in *thousands of dollars*, they’d be represented by the triplets (25,54,700) and (35,76,730), which would give all three variables similar kinds of influence.

Ultimately the way you scale your variables, or equivalently in this situation the way you define your concept of distance, has a potentially enormous effect on the output. In statistics it is called your “prior.”

Euclidean distance is a good go-to distance metric for attributes that are real-valued and can be plotted on a plane or in multidimensional space. Some others are:

Cosine Similarity

Also can be used between two real-valued vectors, \vec{x} and \vec{y} , and will yield a value between -1 (exact opposite) and 1 (exactly the same) with 0 in between meaning independent. Recall the definition $\cos(\vec{x}, \vec{y}) = \frac{\vec{x} \cdot \vec{y}}{\|\vec{x}\| \|\vec{y}\|}$.

Jaccard Distance or Similarity

This gives the distance between a set of objects—for example, a list of Cathy’s friends $A = \{Kahn, Mark, Laura, \dots\}$ and a list of Rachel’s friends $B = \{Mladen, Kahn, Mark, \dots\}$ —and says how similar those two sets are: $J(A, B) = \frac{|A \cap B|}{|A \cup B|}$.

Mahalanobis Distance

Also can be used between two real-valued vectors and has the advantage over Euclidean distance that it takes into account correlation and is scale-invariant. $d(\vec{x}, \vec{y}) = \sqrt{(\vec{x} - \vec{y})^T S^{-1} (\vec{x} - \vec{y})}$, where S is the covariance matrix.

Hamming Distance

Can be used to find the distance between two strings or pairs of words or DNA sequences of the same length. The distance between olive and ocean is 4 because aside from the “o” the other 4 letters are different. The distance between shoe and hose is 3 because aside from the “e” the other 3 letters are different. You just go through each position and check whether the letters the same in that position, and if not, increment your count by 1.

Manhattan

This is also a distance between two real-valued k -dimensional vectors. The image to have in mind is that of a taxi having to travel the city streets of Manhattan, which is laid out in a grid-like fashion (you can’t cut diagonally across buildings). The distance is therefore defined as $d(\vec{x}, \vec{y}) = \sum_i^k |x_i - y_i|$, where i is the i th element of each of the vectors.

There are many more distance metrics available to you depending on your type of data. We start with a Google search when we’re not sure where to start.

What if your attributes are a mixture of kinds of data? This happens in the case of the movie ratings example: some were numerical attributes, such as budget and number of actors, and one was categorical, genre. But you can always define your own custom distance metric.

For example, you can say if movies are the same genre, that will contribute “0” to their distance. But if they’re of a different genre, that will contribute “10,” where you picked the value 10 based on the fact that this was on the same scale as budget (millions of dollars), which is in

the range of 0 and 100. You could do the same with number of actors. You could play around with the 10; maybe 50 is better.

You'll want to justify why you're making these choices. The justification could be that you tried different values and when you tested the algorithm, this gave the best evaluation metric. Essentially this 10 is either a second tuning parameter that you've introduced into the algorithm on top of the k , or a prior you've put on the model, depending on your point of view and how it's used.

Training and test sets

For any machine learning algorithm, the general approach is to have a training phase, during which you create a model and “train it”; and then you have a testing phase, where you use new data to test how good the model is.

For k -NN, the training phase is straightforward: it's just reading in your data with the “high” or “low” credit data points marked. In testing, you pretend you don't know the true label and see how good you are at guessing using the k -NN algorithm.

To do this, you'll need to save some clean data from the overall data for the testing phase. Usually you want to save randomly selected data, let's say 20%.

Your R console might look like this:

```
> head(data)
  age income credit
1  69     3    low
2  66    57    low
3  49    79    low
4  49    17    low
5  58    26   high
6  44    71   high

n.points <- 1000 # number of rows in the dataset
sampling.rate <- 0.8

# we need the number of points in the test set to calculate
# the misclassification rate
num.test.set.labels <- n.points * (1 - sampling.rate)

# randomly sample which rows will go in the training set
training <- sample(1:n.points, sampling.rate * n.points,
                  replace=FALSE)
train <- subset(data[training, ], select = c(Age, Income))
# define the training set to be those rows
```

```
# the other rows are going into the test set
testing <- setdiff(1:n.points, training)
# define the test set to be the other rows
test <- subset(data[testing, ], select = c(Age, Income))

cl <- data$Credit[training]
# this is the subset of labels for the training set
true.labels <- data$Credit[testing]
# subset of labels for the test set, we're withholding these
```

Pick an evaluation metric

How do you evaluate whether your model did a good job?

This isn't easy or universal—you may decide you want to penalize certain kinds of misclassification more than others. False negatives may be way worse than false positives. Coming up with the evaluation metric could be something you work on with a domain expert.

For example, if you were using a classification algorithm to predict whether someone had cancer or not, you would want to minimize false negatives (misdiagnosing someone as not having cancer when they actually do), so you could work with a doctor to tune your evaluation metric.

Note you want to be careful because if you really wanted to have *no* false negatives, you could just tell *everyone* they have cancer. So it's a trade-off between *sensitivity* and *specificity*, where sensitivity is here defined as the probability of correctly diagnosing an ill patient as ill; specificity is here defined as the probability of correctly diagnosing a well patient as well.



Other Terms for Sensitivity and Specificity

Sensitivity is also called the *true positive rate* or *recall* and varies based on what academic field you come from, but they all mean the same thing. And *specificity* is also called the *true negative rate*. There is also the *false positive rate* and the *false negative rate*, and these don't get other special names.

Another evaluation metric you could use is *precision*, defined in [Chapter 5](#). The fact that some of the same formulas have different names is due to the fact that different academic disciplines have developed these ideas separately. So *precision* and *recall* are the quantities

used in the field of information retrieval. Note that *precision* is not the same thing as *specificity*.

Finally, we have *accuracy*, which is the ratio of the number of correct labels to the total number of labels, and the misclassification rate, which is just $1 - \text{accuracy}$. Minimizing the *misclassification rate* then just amounts to maximizing *accuracy*.

Putting it all together

Now that you have a distance measure and an evaluation metric, you're ready to roll.

For each person in your test set, you'll pretend you don't know his label. Look at the labels of his three nearest neighbors, say, and use the label of the majority vote to label him. You'll label all the members of the test set and then use the misclassification rate to see how well you did. All this is done automatically in R, with just this single line of R code:

```
knn (train, test, cl, k=3)
```

Choosing k

How do you choose k ? This is a parameter you have control over. You might need to understand your data pretty well to get a good guess, and then you can try a few different k 's and see how your evaluation changes. So you'll run k-nn a few times, changing k , and checking the evaluation metric each time.



Binary Classes

When you have binary classes like “high credit” or “low credit,” picking k to be an odd number can be a good idea because there will always be a majority vote, no ties. If there is a tie, the algorithm just randomly picks.

```
# we'll loop through and see what the misclassification rate
# is for different values of k
for (k in 1:20) {
  print(k)
  predicted.labels <- knn(train, test, cl, k)
  # We're using the R function knn()
  num.incorrect.labels <- sum(predicted.labels != true.labels)
  misclassification.rate <- num.incorrect.labels /
    num.test.set.labels
}
```

```
    print(misclassification.rate)
  }
```

Here's the output in the form (k, misclassification rate):

```
k  misclassification.rate
1, 0.28
2, 0.315
3, 0.26
4, 0.255
5, 0.23
6, 0.26
7, 0.25
8, 0.25
9, 0.235
10, 0.24
```

So let's go with $k=5$ because it has the lowest misclassification rate, and now you can apply it to your guy who is 57 with a \$37,000 salary. In the R console, it looks like:

```
> test <- c(57,37)
> knn(train,test,cl, k = 5)
[1] low
```

The output by majority vote is a low credit score when $k = 5$.



Test Set in k-NN

Notice we used the function `knn()` twice and used it in different ways. In the first way, the test set was some data we were using to evaluate how good the model was. In the second way, the “test” set was actually a new data point that we wanted a prediction for. We could also have given it many rows of people who we wanted predictions for. But notice that R doesn't know the difference whether what you're putting in for the test set is truly a “test” set where you know the real labels, or a test set where you don't know and want predictions.

What are the modeling assumptions?

In the previous chapter we discussed modeling and modeling assumptions. So what were the modeling assumptions here?

The k-NN algorithm is an example of a nonparametric approach. You had no modeling assumptions about the underlying data-generating distributions, and you weren't attempting to estimate any parameters. But you still made *some* assumptions, which were:

- Data is in some feature space where a notion of “distance” makes sense.
- Training data has been labeled or classified into two or more classes.
- You pick the number of neighbors to use, k .
- You’re assuming that the *observed features* and the *labels* are somehow associated. They may not be, but ultimately your evaluation metric will help you determine how good the algorithm is at labeling. You might want to add more features and check how that alters the evaluation metric. You’d then be tuning both *which* features you were using and k . But as always, you’re in danger here of overfitting.

Both linear regression and k-NN are examples of “supervised learning,” where you’ve observed both x and y , and you want to know the function that brings x to y . Next up, we’ll look at an algorithm you can use when you don’t know what the right answer is.

k-means

So far we’ve only seen supervised learning, where we know beforehand what label (aka the “right answer”) is and we’re trying to get our model to be as accurate as possible, defined by our chosen evaluation metric.

k-means is the first *unsupervised* learning technique we’ll look into, where the goal of the algorithm is to determine the definition of the right answer by finding clusters of data for you.

Let’s say you have some kind of data at the user level, e.g., Google+ data, survey data, medical data, or SAT scores.

Start by adding structure to your data. Namely, assume each row of your dataset corresponds to a user as follows:

```
age gender income state household size
```

Your goal is to *segment* the users. This process is known by various names: besides being called segmenting, you could say that you’re going to *stratify*, *group*, or *cluster* the data. They all mean finding similar types of users and bunching them together.

Why would you want to do this? Here are a few examples:

- You might want to give different users different experiences. Marketing often does this; for example, to offer toner to people who are known to own printers.
- You might have a model that works better for specific groups. Or you might have different models for different groups.
- **Hierarchical modeling** in statistics does something like this; for example, to separately model geographical effects from household effects in survey results.

To see why an algorithm like this might be useful, let's first try to construct something by hand. That might mean you'd bucket users using handmade thresholds.

So for an attribute like age, you'd create bins: 20–24, 25–30, etc. The same technique could be used for other attributes like income. States or cities are in some sense their own buckets, but you might want fewer buckets, depending on your model and the number of data points. In that case, you could bucket the buckets and think of “East Coast” and “Midwest” or something like that.

Say you've done that for each attribute. You may have 10 age buckets, 2 gender buckets, and so on, which would result in $10 \times 2 \times 50 \times 10 \times 3 = 30,000$ possible bins, which is big.

Imagine this data existing in a five-dimensional space where each axis corresponds to one attribute. So there's a gender axis, an income axis, and so on. You can also label the various possible buckets along the corresponding axes, and if you did so, the resulting grid would consist of every possible bin—a bin for each possible combination of attributes.

Each user would then live in one of those 30,000 five-dimensional cells. But wait, it's highly unlikely you'd want to build a different marketing campaign for each bin. So you'd have to bin the bins...

Now you likely see the utility of having an algorithm to do this for you, especially if you could choose beforehand how many bins you want. That's exactly what k-means is: a *clustering* algorithm where *k* is the number of bins.

2D version

Let's back up to a simpler example than the five-dimensional one we just discussed. Let's say you have users where you know how many ads

have been shown to each user (the number of impressions) and how many times each has clicked on an ad (number of clicks).

Figure 3-9 shows a simplistic picture that illustrates what this might look like.

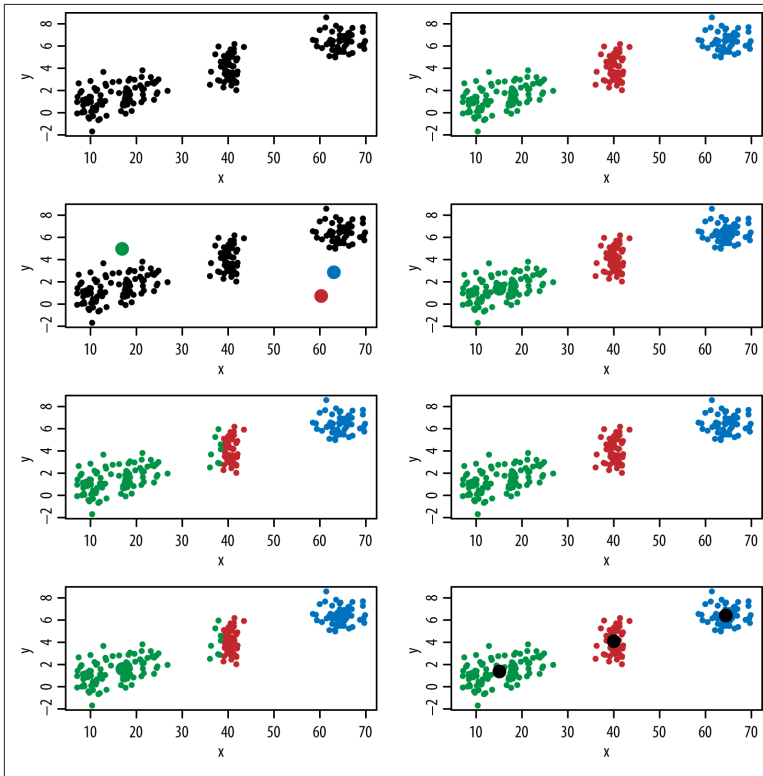


Figure 3-9. Clustering in two dimensions; look at the panels in the left column from top to bottom, and then the right column from top to bottom

Visually you can see in the top-left that the data naturally falls into clusters. This may be easy for you to do with your eyes when it's only in two dimensions and there aren't that many points, but when you get to higher dimensions and more data, you need an algorithm to help with this pattern-finding process. k-means algorithm looks for clusters in d dimensions, where d is the number of features for each data point.

Here's how the algorithm illustrated in [Figure 3-9](#) works:

1. Initially, you randomly pick k centroids (or points that will be the center of your clusters) in d -space. Try to make them near the data but different from one another.
2. Then assign each data point to the closest centroid.
3. Move the centroids to the average location of the data points (which correspond to users in this example) assigned to it.
4. Repeat the preceding two steps until the assignments don't change, or change very little.

It's up to you to interpret if there's a natural way to describe these groups once the algorithm's done. Sometimes you'll need to jiggle around k a few times before you get natural groupings.

This is an example of *unsupervised* learning because the labels are not known and are instead discovered by the algorithm.

k-means has some known issues:

- Choosing k is more an art than a science, although there are bounds: $1 \leq k \leq n$, where n is number of data points.
- There are convergence issues—the solution can fail to exist, if the algorithm falls into a loop, for example, and keeps going back and forth between two possible solutions, or in other words, there isn't a single unique solution.
- Interpretability can be a problem—sometimes the answer isn't at all useful. Indeed that's often the biggest problem.

In spite of these issues, it's pretty fast (compared to other clustering algorithms), and there are broad applications in marketing, computer vision (partitioning an image), or as a starting point for other models.

In practice, this is just one line of code in R:

```
kmeans(x, centers, iter.max = 10, nstart = 1,  
       algorithm = c("Hartigan-Wong", "Lloyd", "Forgy",  
                    "MacQueen"))
```

Your dataset needs to be a matrix, x , each column of which is one of your features. You specify k by selecting centers. It defaults to a certain number of iterations, which is an argument you can change. You can also select the specific algorithm it uses to discover the clusters.

Historical Perspective: k-means

Wait, didn't we just describe the algorithm? It turns out there's more than one way to go after k-means clustering.

The standard k-means algorithm is attributed to separate work by Hugo Steinhaus and Stuart Lloyd in 1957, but it wasn't called "k-means" then. The first person to use that term was James MacQueen in 1967. It wasn't published outside Bell Labs until 1982.

Newer versions of the algorithm are Hartigan-Wong and Lloyd and Forgy, named for their inventors and developed throughout the '60s and '70s. The algorithm we described is the default, Hartigan-Wong. It's fine to use the default.

As history keeps marching on, it's worth checking out the more recent k-means++ developed in 2007 by David Arthur and Sergei Vassilvitskii (now at Google), which helps avoid convergence issues with k-means by optimizing the initial seeds.

Exercise: Basic Machine Learning Algorithms

Continue with the NYC (Manhattan) Housing dataset you worked with in the preceding chapter: <http://abt.cm/1g3A12P>.

- Analyze sales using regression with any predictors you feel are relevant. Justify why regression was appropriate to use.
- Visualize the coefficients and fitted model.
- Predict the neighborhood using a k-NN classifier. Be sure to withhold a subset of the data for testing. Find the variables and the k that give you the lowest prediction error.
- Report and visualize your findings.
- Describe any decisions that could be made or actions that could be taken from this analysis.

Solutions

In the preceding chapter, we showed how explore and clean this dataset, so you'll want to do that first before you build your regression model. Following are two pieces of R code. The first shows how you

might go about building your regression models, and the second shows how you might clean and prepare your data and then build a k-NN classifier.

Sample R code: Linear regression on the housing dataset

Author: Ben Reddy

```
model1 <- lm(log(sale.price.n) ~ log(gross.sqft), data=bk.homes)
## what's going on here?

bk.homes[which(bk.homes$gross.sqft==0),]

bk.homes <- bk.homes[which(bk.homes$gross.sqft>0 &
  bk.homes$land.sqft>0),]
model1 <- lm(log(sale.price.n) ~ log(gross.sqft), data=bk.homes)
summary(model1)

plot(log(bk.homes$gross.sqft), log(bk.homes$sale.price.n))
abline(model1, col="red", lwd=2)
plot(resid(model1))

model2 <- lm(log(sale.price.n) ~ log(gross.sqft) +
  log(land.sqft) + factor(neighborhood), data=bk.homes)
summary(model2)
plot(resid(model2))

## leave out intercept for ease of interpretability
model2a <- lm(log(sale.price.n) ~ 0 + log(gross.sqft) +
  log(land.sqft) + factor(neighborhood), data=bk.homes)
summary(model2a)
plot(resid(model2a))

## add building type
model3 <- lm(log(sale.price.n) ~ log(gross.sqft) +
  log(land.sqft) + factor(neighborhood) +
  factor(building.class.category), data=bk.homes)
summary(model3)
plot(resid(model3))

## interact neighborhood and building type
model4 <- lm(log(sale.price.n) ~ log(gross.sqft) +
  log(land.sqft) + factor(neighborhood)*
  factor(building.class.category), data=bk.homes)
summary(model4)
plot(resid(model4))
```

Sample R code: K-NN on the housing dataset

Author: Ben Reddy
require(gdata)
require(geoPlot)

```

require(class)

setwd("~/Documents/Teaching/Stat 4242 Fall 2012/Homework 2")

mt <- read.xls("rollingsales_manhattan.xls",
  pattern="BOROUGH",stringsAsFactors=FALSE)
head(mt)
summary(mt)

names(mt) <- tolower(names(mt))

mt$sale.price.n <- as.numeric(gsub("[^[:digit:]]", "",
  mt$sale.price))

sum(is.na(mt$sale.price.n))
sum(mt$sale.price.n==0)

names(mt) <- tolower(names(mt))

## clean/format the data with regular expressions
mt$gross.sqft <- as.numeric(gsub("[^[:digit:]]", "",
  mt$gross.square.feet))
mt$land.sqft <- as.numeric(gsub("[^[:digit:]]", "",
  mt$land.square.feet))

mt$sale.date <- as.Date(mt$sale.date)
mt$year.built <- as.numeric(as.character(mt$year.built))
mt$zip.code <- as.character(mt$zip.code)

## - standardize data (set year built start to 0; land and
gross sq ft; sale price (exclude $0 and possibly others); possi
bly tax block; outside dataset for coords of tax block/lot?)
min_price <- 10000
mt <- mt[which(mt$sale.price.n>=min_price),]

n_obs <- dim(mt)[1]

mt$address.noapt <- gsub("[,][[:print:]]*", "",
  gsub("[ ]+", " ", trim(mt$address)))

mt_add <- unique(data.frame(mt$address.noapt,mt$zip.code,
  stringsAsFactors=FALSE))
names(mt_add) <- c("address.noapt", "zip.code")
mt_add <- mt_add[order(mt_add$address.noapt),]

##find duplicate addresses with different zip codes
dup <- duplicated(mt_add$address.noapt)
# remove them
dup_add <- mt_add[mt_add$dup,1]
mt_add <- mt_add[(mt_add$address.noapt != dup_add[1] &
  mt_add$address.noapt != dup_add[2]),]

```

```

n_add <- dim(mt_add)[1]

# sample 500 addresses so we don't run over our Google Maps
API daily limit (and so we're not waiting forever)
n_sample <- 500
add_sample <- mt_add[sample.int(n_add,size=n_sample),]

# first, try a query with the addresses we have
query_list <- addrListLookup(data.frame(1:n_sample,
  add_sample$address.noapt,rep("NEW YORK",times=n_sample),
  rep("NY",times=n_sample),add_sample$zip.code,
  rep("US",times=n_sample))),[,1:4]

query_list$matched <- (query_list$latitude != 0)

unmatched_inds <- which(!query_list$matched)
unmatched <- length(unmatched_inds)

# try changing EAST/WEST to E/W
query_list[unmatched_inds,1:4] <- addrListLookup
  (data.frame(1:unmatched,gsub(" WEST ", " W ",
  gsub(" EAST ", " E ",add_sample[unmatched_inds,1])),
  rep("NEW YORK",times=unmatched), rep("NY",times=unmatched),
  add_sample[unmatched_inds,2],rep("US",times=unmatched))),
1:4]

query_list$matched <- (query_list$latitude != 0)
unmatched_inds <- which(!query_list$matched)
unmatched <- length(unmatched_inds)

# try changing STREET/AVENUE to ST/AVE
query_list[unmatched_inds,1:4] <- addrListLookup
  (data.frame(1:unmatched,gsub(" WEST ", " W ",
  gsub(" EAST ", " E ",gsub(" STREET", " ST",
  gsub(" AVENUE", " AVE",add_sample[unmatched_inds,1]))),
  rep("NEW YORK",times=unmatched), rep("NY",times=unmatched),
  add_sample[unmatched_inds,2],rep("US",times=unmatched))),
1:4]

query_list$matched <- (query_list$latitude != 0)
unmatched_inds <- which(!query_list$matched)
unmatched <- length(unmatched_inds)

## have to be satisfied for now
add_sample <- cbind(add_sample,query_list$latitude,
  query_list$longitude)
names(add_sample)[3:4] <- c("latitude", "longitude")

add_sample <- add_sample[add_sample$latitude!=0,]

add_use <- merge(mt,add_sample)

```

```

add_use <- add_use[!is.na(add_use$latitude),]

# map coordinates
map_coords <- add_use[,c(2,4,26,27)]
table(map_coords$neighborhood)
map_coords$neighborhood <- as.factor(map_coords$neighborhood)

geoPlot(map_coords, zoom=12, color=map_coords$neighborhood)

## - knn function
## - there are more efficient ways of doing this,
## but oh well...

map_coords$class <- as.numeric(map_coords$neighborhood)
n_cases <- dim(map_coords)[1]
split <- 0.8

train_inds <- sample.int(n_cases, floor(split*n_cases))
test_inds <- (1:n_cases)[-train_inds]

k_max <- 10
knn_pred <- matrix(NA, ncol=k_max, nrow=length(test_inds))
knn_test_error <- rep(NA, times=k_max)

for (i in 1:k_max) {
  knn_pred[,i] <- knn(map_coords[train_inds,3:4],
    map_coords[test_inds,3:4], cl=map_coords[train_inds,5], k=i)
  knn_test_error[i] <- sum(knn_pred[,i] !=
    map_coords[test_inds,5])/length(test_inds)
}

plot(1:k_max, knn_test_error)

```

Modeling and Algorithms at Scale

The data you've been dealing with so far in this chapter has been pretty small on the Big Data spectrum. What happens to these models and algorithms when you have to scale up to massive datasets?

In some cases, it's entirely appropriate to sample and work with a smaller dataset, or to run the same model across multiple *sharded* datasets. (Sharding is where the data is broken up into pieces and divided among different machines, and then you look at the empirical distribution of the estimators across models.) In other words, there are statistical solutions to these engineering challenges.

However, in some cases we want to fit these models at scale, and the challenge of scaling up models generally translates to the challenge of

creating parallelized versions or approximations of the optimization methods. Linear regression at scale, for example, relies on matrix inversions or approximations of matrix inversions.

Optimization with Big Data calls for new approaches and theory—this is the frontier! From a 2013 talk by Peter Richtarik from the University of Edinburgh: “In the Big Data domain classical approaches that rely on optimization methods with multiple iterations are not applicable as the computational cost of even a single iteration is often too excessive; these methods were developed in the past when problems of huge sizes were rare to find. We thus need new methods which would be simple, gentle with data handling and memory requirements, and scalable. Our ability to solve truly huge scale problems goes hand in hand with our ability to utilize modern parallel computing architectures such as multicore processors, graphical processing units, and computer clusters.”

Much of this is outside the scope of the book, but a data scientist needs to be aware of these issues, and some of this is discussed in [Chapter 14](#).

Summing It All Up

We’ve now introduced you to three algorithms that are the basis for the solutions to many real-world problems. If you understand these three, you’re already in good shape. If you don’t, don’t worry, it takes a while to sink in.

Regression is the basis of many forecasting and classification or prediction models in a variety of contexts. We showed you how you can predict a continuous outcome variable with one or more predictors. We’ll revisit it again in [Chapter 5](#), where we’ll learn *logistic* regression, which can be used for classification of binary outcomes; and in [Chapter 6](#), where we see it in the context of time series modeling. We’ll also build up your feature selection skills in [Chapter 7](#).

k-NN and k-means are two examples of clustering algorithms, where we want to group together similar objects. Here the notions of *distance* and *evaluation measures* became important, and we saw there is some subjectivity involved in picking these. We’ll explore clustering algorithms including Naive Bayes in the next chapter, and in the context of social networks ([Chapter 10](#)). As we’ll see, *graph clustering* is an interesting area of research. Other examples of clustering algorithms

not explored in this book are *hierarchical clustering* and *model-based clustering*.

For further reading and a more advanced treatment of this material, we recommend the standard classic Hastie and Tibshirani book, *Elements of Statistical Learning* (Springer). For an in-depth exploration of building regression models in a Bayesian context, we highly recommend Andrew Gelman and Jennifer Hill's *Data Analysis using Regression and Multilevel/Hierarchical Models*.

Thought Experiment: Automated Statistician

Rachel attended a workshop in Big Data Mining at Imperial College London in May 2013. One of the speakers, Professor Zoubin Ghahramani from Cambridge University, said that one of his long-term research projects was to build an “automated statistician.” What do you think that means? What do you think would go into building one?

Does the idea scare you? Should it?